

Listing of the claims

1. (Previously Presented) A program storage device readable by a machine, tangibly embodying a program of instructions executable by the machine to perform method steps for rescoring N-best hypotheses of an automatic speech recognition system, the method steps comprising:

generating a synthetic waveform for each of N text sequences representing the N-best hypotheses output from a speech recognition system;

comparing each synthetic waveform with an original waveform decoded by the speech recognition system to determine the synthetic waveform that is closest to the original waveform; and

selecting for output the text sequence corresponding to the synthetic waveform determined to be closest to the original waveform

2. (Original) The program storage device of claim 1, wherein the instructions for performing the comparing step include instructions for performing the steps of:

aligning frames of the original waveform and frames of each synthetic waveform to a corresponding one of the N text sequences; and

calculating the distance between the original waveform and each of the synthetic waveforms based on the corresponding alignments.

3. (Original) The program storage device of claim 2, wherein the instructions for performing the comparing step further include instructions for:

retrieving feature vectors corresponding to the original waveform; and

generating feature vectors for each synthetic waveform such that the feature vectors for the synthetic waveforms are similar in structure to the feature vectors of the original waveform;

wherein the alignment is performed by time-aligning the feature vectors of the original waveform and the feature vectors of each synthetic waveform with the corresponding one of the N text sequences.

4. (Original) The program storage device of claim 2, wherein the alignment is performed using Viterbi alignment process.

5. (Original) The program storage device of claim 2, wherein the alignment is performed on a phoneme level.

6. (Original) The program storage device of claim 2, wherein the instructions for calculating the distance include instructions for performing the steps of:

calculating an individual distance between each aligned frame of the original waveform and each of the N synthetic waveforms; and

summing the individual distances of the aligned frames of the original waveform and each synthetic waveform.

7. (Original) The program storage device of claim 1, wherein the instructions for performing the comparing step include instructions for performing the steps of:

(a) setting a parameter $N=1$;

(b) retrieving the Nth synthetic waveform and the corresponding Nth text sequence;

- (c) time-aligning frames of the original waveform and frames of the Nth synthetic waveform to corresponding text of the Nth text sequence;
- (d) computing an individual distance between each corresponding aligned frame of the original and Nth synthetic waveform;
- (e) summing the individual distances to compute the distance between the original waveform and the Nth synthetic waveform;
- (f) determining if the computed distance is less than a current best distance value;
- (g) setting the current best distance value equal to the computed distance and saving the Nth text sequence for consideration as the final output, if the computed distance is determined to be less than the current best distance threshold;
- (h) incrementing the parameter N by one; and
- (i) repeating steps (b) through (h) until each of the N text sequences have been considered.

8. (Original) The program storage device of claim 7, wherein the instructions for performing the step of determining the individual distance (step d) include instructions for:
computing a mean feature vector of all feature vectors comprising each aligned frame for both the original and Nth synthetic waveform, wherein the individual distance for each aligned frame is calculated by determining a distance between each mean of the corresponding aligned frames.

9. (Previously Presented) A method for rescoring N-best hypotheses of an automatic speech recognition system, the method comprising the steps of:

generating a synthetic waveform for each of N text sequences representing the N-best hypotheses output from a speech recognition system;

comparing each synthetic waveform with an original waveform decoded by the speech recognition system to determine the synthetic waveform that is closest to the original waveform; and

selecting for output the text sequence corresponding to the synthetic waveform determined to be closest to the original waveform.

10. (Original) The method of claim 9, wherein the comparing step includes the steps of:

aligning frames of the original waveform and frames of each synthetic waveform to a corresponding one of the N text sequences; and

calculating the distance between the original waveform and each of the synthetic waveforms based on the corresponding-alignments.

11. (Original) The method of claim 10, wherein the comparing step further includes the steps of:

retrieving feature vectors corresponding to the original waveform; and
generating feature vectors for each synthetic waveform such that the feature vectors for the synthetic waveforms are-similar in structure to the feature vectors of the original waveform;

wherein the alignment is performed by
time-aligning the feature vectors of the original waveform and the feature vectors of each synthetic waveform with the corresponding one of the N text sequences.

12. (Original) The method of claim 10, wherein the step of calculating the distance includes the steps of:

calculating an individual distance between each aligned frame of the original waveform and each of the N synthetic waveforms; and summing the individual distances of the aligned frames of the original waveform and each synthetic waveform.

13. (Original) The method of claim 9, wherein the comparing step includes the steps of:

- (a) setting a parameter $N=1$;
- (b) retrieving the Nth synthetic waveform and the corresponding Nth text sequence;
- (c) time-aligning frames of the original waveform and frames of the Nth synthetic waveform to corresponding text of the Nth text sequence;
- (d) computing an individual distance between each corresponding aligned frame of the original and Nth synthetic waveform;
- (e) summing the individual distances to compute the distance between the original waveform and the Nth synthetic waveform;
- (f) determining if the computed distance is less than a current best distance value;
- (g) setting the current best distance value equal to the computed distance and saving the Nth text sequence for consideration as the final output, if the computed distance is determined to be less than the current best distance threshold;
- (h) incrementing the parameter N by one; and

(i) repeating steps (b) through (h) until each of the N text sequences have been considered.

14. (Original) The method of claim 13, wherein the step of determining the individual distance (step d) includes the steps of:

computing a mean feature vector of all feature vectors comprising each aligned frame for both the original and Nth synthetic waveform, wherein the individual distance for each aligned frame is calculated by determining a distance between each means of the corresponding aligned frames.

15. (Original) An automatic speech recognition system, comprising:

a decoder for decoding an original waveform of acoustic utterances to produce N text sequences, the N text sequences representing N-best hypotheses of the decoded original waveform;

a waveform generator for generating a synthetic waveform for each of the N text sequences; and

a comparator for comparing each synthetic waveform with the original waveform to rescore the N-best hypotheses.

16. (Original) The system of claim 15, further comprising a feature analysis processor adapted to generate a set of feature vectors for the original waveform and generate a set of feature vectors for each of the N synthetic waveforms using a similar feature analysis process.

17. (Original) The system of claim 15, further comprising a processor adapted to process one of the original waveform, the synthetic-waveforms, and both, to compensate for speaker-dependent variations.

18. (Original) The system of claim 15, wherein the comparator comprises:
means for determining the synthetic waveform that is closest in distance to the original waveform; and
means for outputting the N text sequence corresponding to the synthetic waveform that is determined to be closest to the original waveform.

19. (Original) The system of claim 18, wherein the means for determining the closest synthetic waveform utilizes one of a distance score, a language model score, an acoustic model score, and a combination thereof, for determining the closest distance.

20. (Original) The system of claim 18, wherein the means for determining the closest synthetic waveform comprises:
means for aligning frames of the original waveform and frames of each synthetic waveform to a corresponding one of the N text sequences; and
means for calculating the distance between the original waveform and each of the synthetic waveforms based on the corresponding alignments.

21. (Original) The system of claim 20, wherein the frames are aligned on a phoneme level.

22. (Original) The system of claim 20, wherein the means for calculating the distance comprises:

means for calculating an individual distance between each aligned frame of the original waveform and each of the N synthetic waveforms; and

means for summing the individual distances of the aligned frames of the original waveform and each synthetic waveform to compute the distance between the original waveform and each synthetic waveforms.